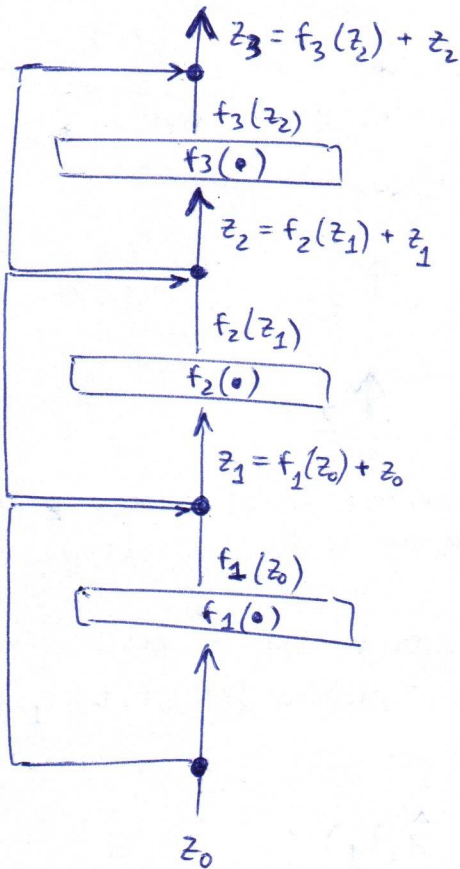


VIDEO "Neural Ordinary Differential Equations"
- by "Andriy Drozdyuk"

ResNET



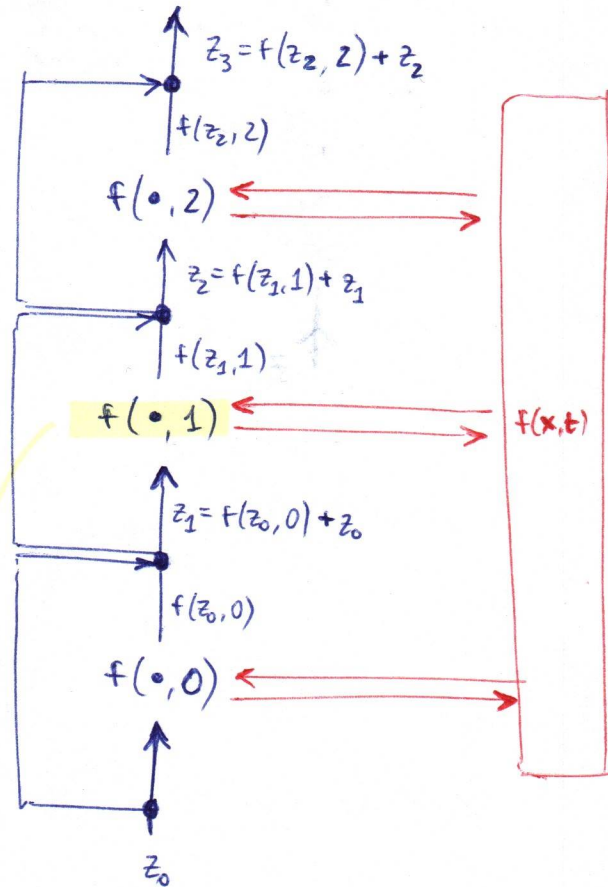
$$z_{i+1} = \underbrace{f_{i+1}(z_i)}_{\text{residual}} + \underbrace{z_i}_{\text{skip connection}}$$

the output of every layer is something done to the input ($f_i(z_i)$) + the original input (z_i).

(Note: f_i can be anything and $f_i =$ or $\neq f_j$ for $i \neq j$)

VS.

ODE Net



what time does this function visit $f(x, t)$? $t = 1$
who brings with him? $\bullet = z_1$

We can see z as a function of time: $z(t)$:

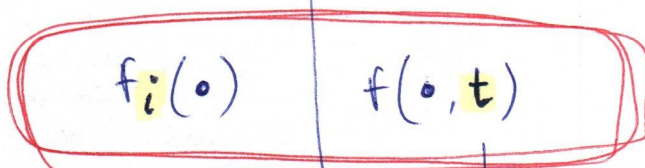
$$z_0 = z(0)$$

$$z_1 = z(1)$$

$$z_2 = z(2)$$

$$z_3 = z(3)$$

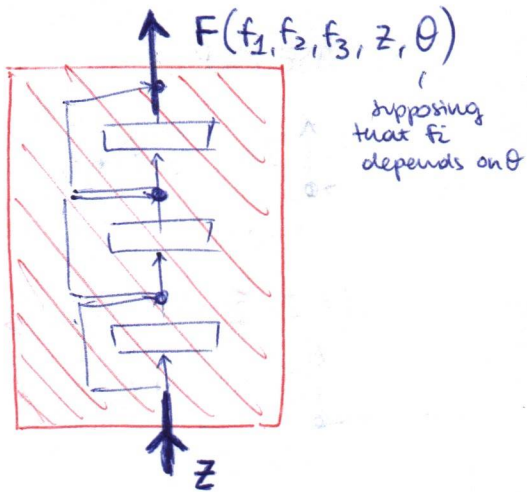
however, we don't know z , it's hidden



dependence based on which moment we use it

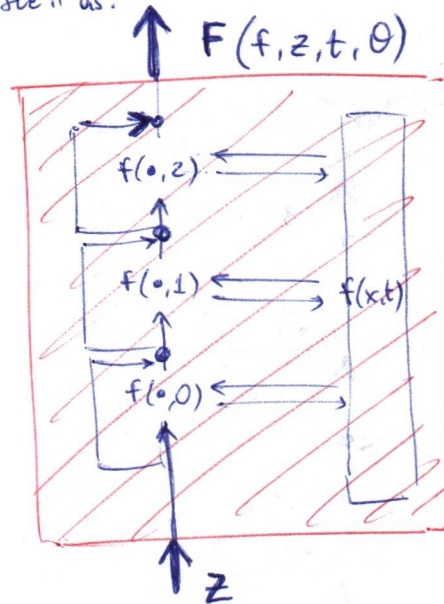
(ResNet)

We can see it as:



(ODENet)

We can see it as:



If we know the initial value of z ($z(0)$) we can move to the field of dynamical systems:

⇒ through ODE solvers
(e.g. $\text{ODESolve}(z(0), f, t_0, t_1, \theta)$)
we get:

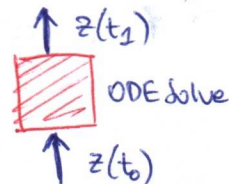
$$\hat{z}(t_1)$$

⇒ the output of the neural network will be the (approximated) evaluation of $z(0)$ at t_1 .

We define the loss function as:

$$\hat{z}(t_1) - z(t_1).$$

Basically, through some black box method (ODE solver) we get the output $z(t_1)$ starting from $z(t_0)$



Once we have $\hat{z}(t_1)$ we evaluate the loss. How do we backpropagate the error through \square ?

⇒ ADJOINT METHOD

Dynamical Systems:

characterized by:

- a state that changes in time:

$$z(t)$$

- a rule for how the state changes:

$$\frac{dz}{dt} = f(t, z, p) \quad \text{ODE}$$

time state parameters

usually this is a system of ODEs

The solution to an ODE is a function: $z(t)$.

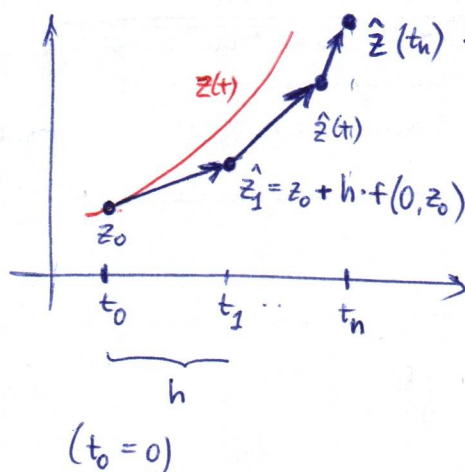
We are usually given with an initial value $z(0)$ (more generally: $z(t_0)$).

Finding a solution to a dynamical system:

Given a dynamical system $f(t, z, p)$ with an initial condition z_0 , find the solution $z(t)$.

Euler method:

$$z_{t+1} = z_t + h f(t, z_t)$$

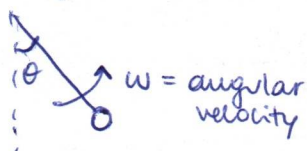


we have solved the system for this time

the divergence of $\hat{z}(t)$ and $z(t)$ is quite a bit!
However it's still something:
we can find a solution $z(t)$ for any t .

Example: Pendulum

(pseudopython)



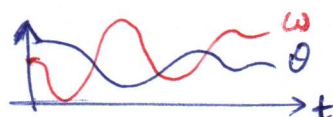
$y_0 = [\pi - 0.1, 0.0]$ — initially pendulum is nearly vertical

$t = \text{np.linspace}(0, 10, 101)$ — 101 points of $0 \leq t \leq 10$

from scipy.integrate import odeint

$\text{sol} = \text{odeint}(\text{pendulum}, y_0, t, \text{args}=(b, c))$

→ $\text{plot sol[:, 0], sol[:, 1]}$:



rule of changing the state

def pendulum(y, t, b, c):

$\theta, w = y$ ← state

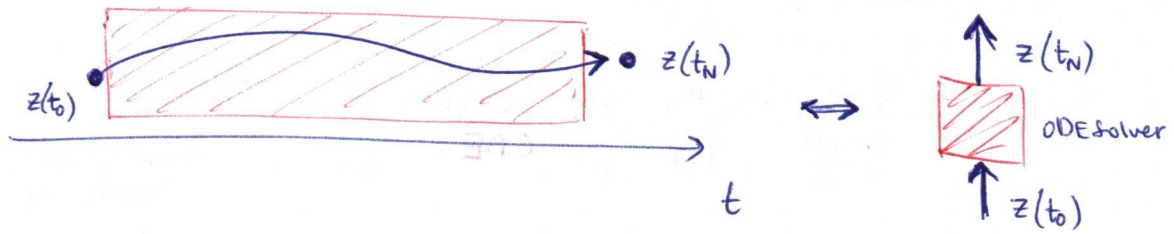
$$\frac{dy}{dt} = [w, -bw - c \sin(\theta)]$$

return $\frac{dy}{dt}$

Adjoint method:

We already know that through an ODE solver we can obtain $z(t_N)$ starting from $z(t_0)$, f (and t_0, t_N, θ)

State space:



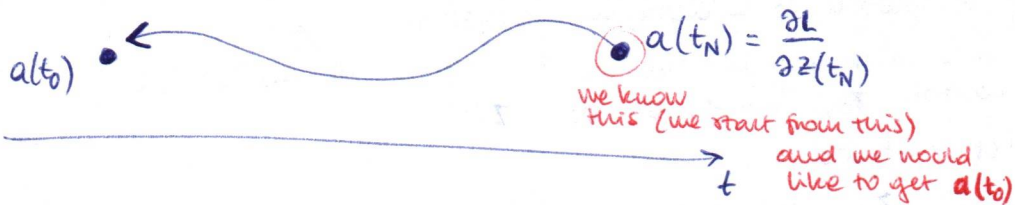
How can we backpropagate the error?

We define an adjoint.

We define:

$$a(t_i) := \frac{\partial L}{\partial z(t_i)}$$

Adjoint state:



→ we can use the same approach as before:

if we treat $a(t_N)$ as initial condition and if we know how $a(\cdot)$ changes w.r.t. time then we can integrate backwards and obtain $a(t_0)$

"continuous-time backpropagation"

ADJOINT METHOD

- Define: $a(t) := \frac{\partial L}{\partial z(t)}$ (previous) black-box ODE solver
- Forward: $z(t+1) = z(t) + \int_t^{t+1} f(z(t)) dt$
- Backward: $a(t) = a(t+1) + \int_{t+1}^t a(t) \frac{\partial f(z(t))}{\partial z(t)} dt$ (1)
adjoint state
- Params: $\frac{\partial L}{\partial \theta_{(2)}} = \int_t^{t+1} a(t) \frac{\partial f(z(t), \theta)}{\partial \theta} dt$

How can we obtain adjoint differential equation (1) and gradients w.r.t. θ (2)?

(Adjoint method - pt. 2)

Let $z(t)$ follow the differential equation: $\frac{dz(t)}{dt} = f(z(t), t, \theta)$, where θ are the parameters.

If we define the adjoint state:

$$a(t) := \frac{\partial L}{\partial z(t)}$$

→ the adjoint state $a(t)$ follows the diff. equation:

$$\frac{da(t)}{dt} = -a(t) \cdot \frac{\partial f(z(t), t, \theta)}{\partial z(t)}$$

adjoint
dynamics

To really conclude we also need to know how θ vary in time.
Let's say that all $[z, \theta, t]$ changes in time.

We call $[z, \theta, t]$ augmented state.

$$\underbrace{\frac{d}{dt} \begin{bmatrix} z \\ \theta \\ t \end{bmatrix}}_{\text{derivative of the augmented state}}(t) = f_{\text{aug}}([z, \theta, t]) := \begin{bmatrix} f([z, t, \theta]) \\ 0 \\ 1 \end{bmatrix}$$

(Hp. $\frac{d\theta}{dt} = 0$)
($\frac{dt}{dt} = 1$)

We define the augmented adjoint as:

$$a_{\text{aug}} = \begin{bmatrix} a \\ a_{\theta} \\ a_t \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial z(t)} \\ \frac{\partial L}{\partial \theta(t)} \\ \frac{\partial L}{\partial t(t)} \end{bmatrix}$$

$$\rightarrow \frac{da_{\text{aug}}(t)}{dt} = - \left[a \frac{\partial f}{\partial z}, \quad a \frac{\partial f}{\partial \theta}, \quad a \frac{\partial f}{\partial t} \right](t)$$

→ **Algorithm 1** (Reverse-mode derivative of an ODE initial value problem)
(paper "Neural Ordinary Differential Equations")